

The upper bound on knots in neural networks

Kevin K. Chen^{*†}

November 2016

Abstract

Neural networks with rectified linear unit activations are essentially multivariate linear splines. As such, one of many ways to measure the “complexity” or “expressivity” of a neural network is to count the number of knots in the spline model. We study the number of knots in fully-connected feedforward neural networks with rectified linear unit activation functions. We intentionally keep the neural networks very simple, so as to make theoretical analyses more approachable. An induction on the number of layers l reveals a tight upper bound on the number of knots in $\mathbb{R} \rightarrow \mathbb{R}^p$ deep neural networks. With $n_i \gg 1$ neurons in layer $i = 1, \dots, l$, the upper bound is approximately $n_1 \dots n_l$. We then show that the exact upper bound is tight, and we demonstrate the upper bound with an example. The purpose of these analyses is to pave a path for understanding the behavior of general $\mathbb{R}^q \rightarrow \mathbb{R}^p$ neural networks.

1 Introduction

In recent years, neural networks—and deep neural networks in particular—have succeeded exceedingly well in such a great plethora of data-driven problems, so as to herald an entire paradigm shift in the way data science is approached. Many everyday computerized tasks—such as image and optical character recognition, the personalization of Internet search results and advertisements, and even playing games such as chess, backgammon, and Go—have been deeply impacted and vastly improved by the application of neural networks. The applications of neural networks, however, have advanced significantly more rapidly than the theoretical understanding of their successes. Elements of neural network structures—such as the division of vector spaces into convex polytopes, and the application of nonlinear activation functions—afford neural networks a great flexibility to model many classes of functions with spectacular accuracy. The flexibility is embodied in universal approximation theorems (Cybenko 1989; Hornik et al. 1989; Hornik 1991; Sonoda and Murata 2015), which essentially state that neural networks can model any continuous function arbitrarily well. The complexity of neural networks, however, have also made their analytical understanding somewhat elusive.

The general thrust of this paper, as well as two companion papers (Chen et al. 2016b,a), is to explore some unsolved elements of neural network theory, and to do so in a way that is independent of specific problems. In the broadest sense, we seek to understand what models neural networks are capable of producing. There exist many variations of neural networks, such as convolutional neural networks, recurrent neural networks, and long short-term memory models, each having their own arenas of success. For simplicity, we choose to focus on the simplest case of feedforward,

^{*}Institute for Defense Analyses, Center for Communications Research - La Jolla

[†]Email for correspondence: kkchen@ccrwest.org

fully-connected neural networks with rectified linear unit activations. This model is defined more precisely in Section 2.

More specifically, as we will see, neural networks with rectified linear unit activations are linear splines; i.e., they are continuous, piecewise linear functions with a finite number of pieces. Therefore, one of many ways to measure of the “complexity” or “expressivity” of a neural network is to count the number of knots, i.e., discontinuities in the first derivative of the output quantities with respect to input quantities. Similarly, one could count the number of piecewise linear regions given by the neural network. Previous works (e.g., Montúfar et al. 2014; Pascanu et al. 2014; Raghu et al. 2016) have observed or shown that number of piecewise linear pieces grows exponentially with the number of layers in the neural network, therefore justifying the use of deep networks over shallow networks.

In this paper, we continue the exploration of how the size of a neural network, given by the *width* or the number of neurons in a layer, and the *depth* or the number of layers, is related to the number of knots in the neural network. Whereas previous works have generally focused on asymptotic or otherwise approximate upper bounds, we derive an exact tight upper bound. The chief utility of such a bound is that it allows an *a priori* determination of whether a neural network size is sufficient for a given task or governing equation. For instance, we could imagine that a neural network designer at least roughly knows the complexity of the input–output behavior of a function to be modeled. In this case, certain neural network widths and depths could be ruled out, on the grounds that no neural networks of those sizes could produce enough knots to model the function of interest.

In this paper, we attempt to circumvent some of the complexities of neural network behavior by making simplifications that may seem strong at times. For instance, the results we report apply specifically to $\mathbb{R} \rightarrow \mathbb{R}^p$ functions. Although neural networks are almost never used to study single-input functions, the simplicity does admit certain analyses that would otherwise be very difficult for general $\mathbb{R}^q \rightarrow \mathbb{R}^p$ functions. Indeed, a key objective following this paper is to extend the results to multidimensional inputs. This extension is tantamount to analyzing convex polytopes in \mathbb{R}^q instead of linear segments in \mathbb{R} in the input space.

The main results of the paper are given by the following theorems.

Theorem 1. *In an l -layer $\mathbb{R} \rightarrow \mathbb{R}^p$ neural network with n_i rectified linear unit neurons in layer $i = 1, \dots, l$, the number of knots m_l in the neural network model satisfies*

$$m_l \leq \sum_{i=1}^l n_i \prod_{j=i+1}^l (n_j + 1). \quad (1)$$

Theorem 2. *If $n_i \geq 3$ for $i = 1, \dots, l - 1$ and $n_l \geq 2$, then the upper bound (1) is tight.*

This paper is organized as follows. Section 2 briefly reviews the neural network architecture that we employ in this paper. Constructive proofs of Theorems 1 and 2 are presented respectively in Sections 3 and 4. An example of a deep neural network meeting the upper bound on the number of knots is then constructed in Section 5. Finally, we summarize our work and comment on future directions in Section 6.

2 Brief overview of neural networks

In Section 2.1, we first review the basic definitions and descriptions of neural networks. Next, we describe two ideas which are relevant for the analytical development of the paper. Section 2.2

describes the rectified linear unit neural network as a linear spline with associated knots and roots, so as to allow knot counting. Afterwards, Section 2.3 derives a transformation of the neural network into an equivalent model with only forward-facing rectified linear units. Such a transformation is useful in constructing particular neural networks (e.g., for Theorem 2 and its associated lemmas).

2.1 Description of neural networks

Neural networks are most commonly employed in the context of supervised machine learning, where the primary objective is to construct a function that best models a data set. In this paper, however, we will be more concerned with the functional behavior of neural network models than with the training of such models. As such, we will not address common topics such as model risk, loss, and optimization. A review of machine learning techniques and their statistical analyses can be found in Knox (2016).

We begin by defining neural networks of a single or multiple hidden layers. It is noteworthy that many variations on neural networks exist. The definitions below correspond to the dense, fully-connected, feedforward structure we will employ, but may differ from architectures used in other studies or applications.

Definition. For some *bias* $b \in \mathbb{R}$, *weight* $\mathbf{w} \in \mathbb{R}^n$, nonlinear *activation function* $\sigma : \mathbb{R} \rightarrow \mathbb{R}$, and *input* $\mathbf{v} \in \mathbb{R}^n$, a *neuron* is the function $\sigma(\mathbf{w} \cdot \mathbf{v} + b)$.

Definition. Let q and p respectively denote the input and output dimension. For $k = 1, \dots, n$, with n the number of neurons, select *input biases* $b_{1k} \in \mathbb{R}$ and *input weights* $\mathbf{w}_{1k} \in \mathbb{R}^q$. Also, for $k = 1, \dots, p$, select *output biases* $b_{2k} \in \mathbb{R}$ and *output weights* $\mathbf{w}_{2k} \in \mathbb{R}^n$. Using the shorthand notation $\mathbf{v} := [v_1 \cdots v_n] \in \mathbb{R}^n$ and $\mathbf{y} := [y_1 \cdots y_p] \in \mathbb{R}^p$, a *single-hidden-layer* neural network is the model $\hat{\mathbf{f}} : \mathbb{R}^q \rightarrow \mathbb{R}^p$, $\mathbf{x} \mapsto \mathbf{y}$ given by

$$v_k := \sigma(\mathbf{w}_{1k} \cdot \mathbf{x} + b_{1k}), \quad k = 1, \dots, n, \quad (2a)$$

$$y_k := \mathbf{w}_{2k} \cdot \mathbf{v} + b_{2k}, \quad k = 1, \dots, p. \quad (2b)$$

This architecture is shown in Figure 1. In summary, each neuron takes an affine transformation of the input and applies the activation function (2a). Then, each output takes an affine transformation of all the neural outputs (2b). The flexibility of this architecture is apparent from the $(q + 1)n + (n + 1)p$ scalars that comprise the biases and weights. In particular, the well-known universal approximation theorem loosely states that if the activation function σ is continuous, non-constant, and bounded, then the single-hidden-layer neural network can approximate any continuous function arbitrarily well with a finite number n of neurons (Cybenko 1989; Hornik et al. 1989; Hornik 1991). A recent result (Sonoda and Murata 2015) extends the universal approximation result to the commonly employed *rectified linear unit*

$$\sigma(x) := \max(0, x) = \frac{x + |x|}{2}. \quad (3)$$

Although the universal approximation theorem implies that the single-hidden-layer neural network is sufficiently flexible for modeling continuous functions, it is common to employ *deep neural networks*, where the outputs of neurons are fed into further hidden layers of neurons. Such architectures are behind many of the notable successes in machine learning applications. The deep neural network with l layers proceeds as follows.

Definition. Let q and p respectively denote the input and output dimension. Set n_i as the number of neurons for each layer $i = 1, \dots, l$. For $k = 1, \dots, n_1$, select input weight vectors $\mathbf{w}_{1k} \in \mathbb{R}^q$ and

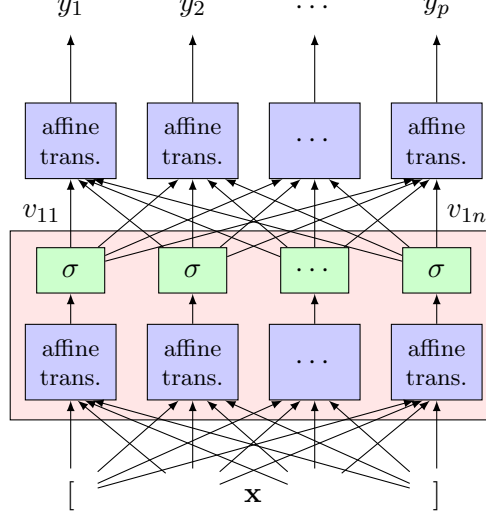


Figure 1: The single-hidden-layer neural network, with the hidden layer shown in red.

input biases $b_{1k} \in \mathbb{R}$. Also, for $i = 2, \dots, l$ and for each $k = 1, \dots, n_i$, also select weight vectors $\mathbf{w}_{ik} \in \mathbb{R}^{n_{i-1}}$ and biases $b_{ik} \in \mathbb{R}$. Finally, for $k = 1, \dots, p$, select output weight vectors $\mathbf{w}_{l+1,k} \in \mathbb{R}^{n_l}$ and output biases $b_{l+1,k} \in \mathbb{R}$. Using the shorthand notation $\mathbf{v}_i := [v_{i1} \cdots v_{in_i}] \in \mathbb{R}^{n_i}$ and $\mathbf{y} := [y_1 \cdots y_p]$, a *deep neural network* is the model $\hat{\mathbf{f}} : \mathbb{R}^q \rightarrow \mathbb{R}^p$, $\mathbf{x} \mapsto \mathbf{y}$ given by

$$v_{1k} := \sigma(\mathbf{w}_{1k} \cdot \mathbf{x} + b_{1k}), \quad k = 1, \dots, n_1 \quad (4a)$$

$$v_{ik} := \sigma(\mathbf{w}_{ik} \cdot \mathbf{v}_{i-1} + b_{ik}), \quad i = 2, \dots, l, \quad k = 1, \dots, n_i \quad (4b)$$

$$y_k := \mathbf{w}_{l+1,k} \cdot \mathbf{v}_l + b_{l+1,k}, \quad k = 1, \dots, p. \quad (4c)$$

The deep neural network architecture is shown in Figure 2. Typically, $n_1 > \cdots > n_l$; it has been empirically shown that training risk is better reduced by optimizing layers closer to the input than layers closer to the output (Raghu et al. 2016).

2.2 Splines, knots, and roots

In this study, we will use the rectified linear unit activation function (3) in all neurons. The rectified linear unit is a common choice because it creates flexible models and is fast to compute. Other common choices, such as the sigmoid function $1/(1 + e^{-x})$, are more computationally intensive. They also typically have smaller regions in the domain where the first derivative is far from zero, which can pose additional challenges when training neural networks on data.

With the rectified linear unit activation, the neural network is essentially a linear spline. To understand this property, first consider the simplified case of a single scalar input, i.e., where the neural network is some $\hat{\mathbf{f}} : \mathbb{R} \rightarrow \mathbb{R}^p$, $x \mapsto \mathbf{y}$. The outputs of the first hidden layer (2a, 4a) are $v_{1k}(x) = \sigma(w_{1k}x + b_{1k})$ for $k = 1, \dots, n_1$. Since $\sigma(x)$ is continuous and has a discontinuity in $d\sigma/dx$ at $x = 0$, v_{1k} is clearly also continuous and has a discontinuity in dv_{1k}/dx at $x = -b_{1k}/w_{1k}$. Thus, the functions $v_{1k}(x)$ are linear splines. The next layer, whether it is a second hidden layer or the output layer, then computes an affine transformation of the functions $v_{1k}(x)$. Such an affine transformation is continuous; hence, it is still a linear spline. This reasoning can be carried out through each hidden layer to the output.

In every application of a rectified linear unit beyond the first layer, knots can be retained, destroyed, or created. An example of this process is shown in Figure 3 for some neuron k in

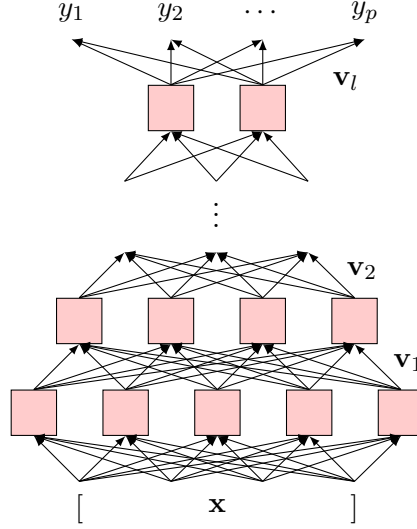


Figure 2: The deep neural network, with each neuron shown in red.

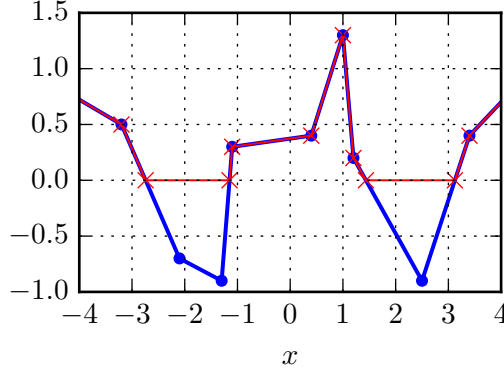


Figure 3: Blue: an example of the affine transformation $\mathbf{w}_{ik} \cdot \mathbf{v}_{i-1}(x) + b_{ik}$ in neuron k of layer i . The knots (filled dots) originate from the various scalar elements of $\mathbf{v}_{i-1}(x)$. Red: the neural output $\sigma(\mathbf{w}_{ik} \cdot \mathbf{v}_{i-1}(x) + b_{ik})$, with knots shown as \times . Knots of the blue spline above zero are retained, knots below zero are discarded, and roots of the blue spline appear as new knots.

some layer i . If the previous layer output $\mathbf{v}_{i-1}(x)$ contains a particular knot at some x_j such that $\mathbf{w}_{ik} \cdot \mathbf{v}_{i-1}(x_j) + b_{ik} > 0$, then the application of the rectified linear unit does not alter this knot, and the knot is retained by this neuron. On the other hand, if $\mathbf{w}_{ik} \cdot \mathbf{v}_{i-1}(x_j) + b_{ik} < 0$, then both the knot and the immediate neighborhood of x_j are rectified to zero, and the knot at x_j is destroyed. Finally, wherever $\mathbf{w}_{ik} \cdot \mathbf{v}_{i-1}(x) + b_{ik}$ crosses zero, there exists a region on one side of the root that is rectified to zero. The rectification introduces a new knot at the root, as shown in Figure 3.

In all three cases, the neural output $\sigma(\mathbf{w}_{ik} \cdot \mathbf{v}_{i-1}(x) + b_{ik})$ again remains a continuous function with discrete discontinuities in its first derivative. Hence, even deep neural networks with rectified linear unit activations are linear splines. The mechanisms for retaining, destroying, and creating knots will be relevant when deriving the upper bound on the number of knots in Section 3. The description of knots becomes more sophisticated in the typical scenario where the input space is \mathbb{R}^q with $q > 1$. In this case, each neuron in the first hidden layer divides the input space into two regions split by the hyperplane $\mathbf{w}_{1k} \cdot \mathbf{x} + b_{1k} = 0$. With the rectified linear unit acting on $\mathbf{w}_{1k} \cdot \mathbf{x} + b_{1k}$, each neuron outputs zero on one side of the hyperplane, and a half-plane with normal vector $[\mathbf{x} \ v_{1k}] =$

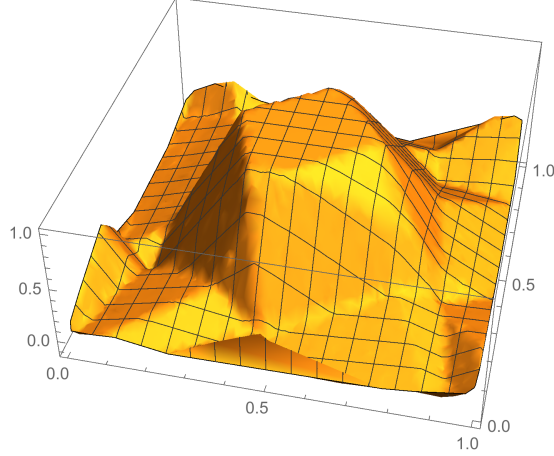


Figure 4: A $\mathbb{R}^2 \rightarrow \mathbb{R}$ neural network model.

$[-\mathbf{w}_{1k} \ 1]$ on the other side. Just as further hidden layers retain, destroy, and create new knots for $q = 1$, further hidden layers retain, destroy, and create new hyperplanes or pieces thereof for $q > 1$. The resulting neural network is a piecewise linear $\mathbb{R}^q \rightarrow \mathbb{R}^p$ model on a finite number of convex polytopes; see Figure 4 for an example. It is still possible to analyze such neural networks in a one-dimensional sense if we were to consider one-dimensional trajectories through the input space \mathbb{R}^q (Raghu et al. 2016), but the full model is notably more complex in general. Many analytical results on multidimensional input spaces rely on upper bounds and asymptotics based on polytope counting (Montúfar et al. 2014; Pascanu et al. 2014; Raghu et al. 2016).

2.3 Equivalent form with forward-facing rectified linear units

For the simple case of $\mathbb{R} \rightarrow \mathbb{R}^p$ neural networks with one hidden layer, the neurons in the first hidden layer (2a, 4a) output $v_{1k} = \sigma(w_{1k}x + b_{1k})$, which is essentially a rectified linear unit $\sigma(x)$ that is horizontally stretched and translated, and possibly reflected across the v_{1k} -axis. Therefore, the sloped ray in the activated region can extend into quadrants I or II in the x - v_{1k} plane. For the purpose of constructing or analyzing $\mathbb{R} \rightarrow \mathbb{R}^p$ neural networks, it is convenient to have all rectified linear units extend in the positive x direction (i.e., into quadrant I), which we call “forward-facing.” Such a feature allows us to consider the action of each rectified linear unit by starting at $x = -\infty$ and increasing x . Thus, no rectified linear units are activated at $x = -\infty$, and the units are successively activated with increasing x ; no units are deactivated.

The transformation that expresses the scalar-input, single-hidden-layer neural network with forward-facing rectified linear units is as follows.

Lemma 1. *Consider the single-hidden-layer $\mathbb{R} \rightarrow \mathbb{R}^p$ rectified linear unit neural network with input weights $w_{1j} \in \mathbb{R}$, input biases $b_{1j} \in \mathbb{R}$, output weights $w_{2kj} \in \mathbb{R}$, and output biases $b_{2k} \in \mathbb{R}$ for $j = 1, \dots, n$ and $k = 1, \dots, p$. The neural network model*

$$y_k(x) = \sum_{j=1}^n w_{2kj} \sigma(w_{1j}x + b_{1j}) + b_{2k}, \quad k = 1, \dots, p \quad (5)$$

(cf. (2) with $\mathbf{w}_{2k} = [w_{2k1} \ \dots \ w_{2kn}]$) is equivalently

$$y_k(x) = \sum_{j=1}^n s_{kj} \sigma(x - x_j) + c_{1k}x + c_{0k}, \quad k = 1, \dots, p, \quad (6)$$

where

$$c_{1k} := \sum_{\substack{1 \leq j \leq n \\ w_{1j} < 0}} w_{2kj} w_{1j}, \quad c_{0k} := \sum_{\substack{1 \leq j \leq n \\ w_{1j} < 0}} w_{2kj} b_{1j} + b_{2k}, \quad (7a)$$

$$s_{kj} := w_{2kj} |w_{1j}|, \quad x_j := -\frac{b_{1j}}{w_{1j}}, \quad j = 1, \dots, n \quad (7b)$$

for $k = 1, \dots, p$. All rectified linear units in (6) face forward.

Proof. We first split the sum in (5) according to the sign of w_{1j} , so that

$$y_k(x) = \sum_{\substack{1 \leq j \leq n \\ w_{1j} < 0}} w_{2kj} \sigma(w_{1j}x + b_{1j}) + \sum_{\substack{1 \leq j \leq n \\ w_{1j} \geq 0}} w_{2kj} \sigma(w_{1j}x + b_{1j}) + b_{2k}. \quad (8)$$

Next, we observe from (3) that

$$\sigma(x) = \sigma(-x) + x; \quad (9)$$

using this property on the first sum, we obtain

$$\begin{aligned} y_k(x) &= \sum_{\substack{1 \leq j \leq n \\ w_{1j} < 0}} w_{2kj} \sigma(-w_{1j}x - b_{1j}) + \sum_{\substack{1 \leq j \leq n \\ w_{1j} < 0}} w_{2kj} (w_{1j}x + b_{1j}) \\ &\quad + \sum_{\substack{1 \leq j \leq n \\ w_{1j} \geq 0}} w_{2kj} \sigma(w_{1j}x + b_{1j}) + b_{2k} \end{aligned} \quad (10a)$$

$$= \sum_{\substack{1 \leq j \leq n \\ w_{1j} < 0}} w_{2kj} \sigma(-w_{1j}x - b_{1j}) + \sum_{\substack{1 \leq j \leq n \\ w_{1j} \geq 0}} w_{2kj} \sigma(w_{1j}x + b_{1j}) + c_{1k}x + c_{0k}. \quad (10b)$$

To combine the two sums, we further observe that if $w \geq 0$, then $\sigma(wx) = w\sigma(x)$. Thus, we can pull $-w_{1j}$ out of the rectified linear unit in the first sum and w_{1j} out of the same in the second sum, and obtain

$$y_k(x) = \sum_{\substack{1 \leq j \leq n \\ w_{1j} < 0}} -w_{2kj} w_{1j} \sigma\left(x + \frac{b_{1j}}{w_{1j}}\right) + \sum_{\substack{1 \leq j \leq n \\ w_{1j} \geq 0}} w_{2kj} w_{1j} \sigma\left(x + \frac{b_{1j}}{w_{1j}}\right) \quad (11a)$$

$$\begin{aligned} &\quad + c_{1k}x + c_{0k} \\ &= \sum_{j=1}^n w_{2kj} |w_{1j}| \sigma\left(x + \frac{b_{1j}}{w_{1j}}\right) + c_{1k}x + c_{0k}, \end{aligned} \quad (11b)$$

which is equal to (6). All rectified linear units face forward because the coefficient on x is simply unity. \square

Besides that all the rectified linear units in (6) face forward, the utility of that expression is that the entire neural network is expressed in terms of four sets of parameters (7), each with a natural interpretation. The parameter x_j is the location of the knot created by neuron j . For convenience, we will assume hereafter that all parameters in j (i.e., w_{1j} , b_{1j} , w_{2kj} , s_{kj} , and x_j) are sorted by ascending x_j . Next, in the contribution from the forward-facing rectified linear unit in neuron j to the scalar output k , s_{kj} is the slope of the activated region. Finally, c_{1k} and c_{0k} describe the line that is added to the sum of rectified linear units, so as to complete the equivalence between (5) and (6).

3 Upper bound on number of knots

Some recent articles have derived asymptotic or otherwise approximate upper bounds for the number of linear regions in neural networks with multidimensional inputs and outputs. For instance, building on Pascanu et al. (2014), Montúfar et al. (2014) showed that for an $\mathbb{R}^q \rightarrow \mathbb{R}^p$ neural network with $n_i \geq q$ neurons in layer $i = 1, \dots, l$, the upper bound on the number of linear regions is at least

$$\left(\prod_{i=1}^{l-1} \left\lfloor \frac{n_i}{q} \right\rfloor^q \right) \sum_{j=0}^q \binom{n_l}{j}. \quad (12)$$

Later, Raghu et al. (2016) gave asymptotic upper bounds for the number of linear regions in neural networks with multidimensional inputs and outputs. The article shows that an $\mathbb{R}^q \rightarrow \mathbb{R}^p$ neural network with n neurons in each of l layers has a number of regions that grows at most like $\mathcal{O}(n^{ql})$ for rectified linear unit activations, and $\mathcal{O}((2n)^{ql})$ for step activation functions. Furthermore, the asymptotic upper bound is shown to be tight (Montúfar et al. 2014; Pascanu et al. 2014; Raghu et al. 2016).

In this section, we derive an exact as opposed to asymptotic or approximate upper bound, but restrict ourselves to the case of $\mathbb{R} \rightarrow \mathbb{R}^p$ neural networks. The possibility of extending the result to $\mathbb{R}^q \rightarrow \mathbb{R}^p$ remains open. We first discuss the mechanisms by which the maximal number of knots is retained and created in each hidden layer. Next, we use induction to prove Theorem 1, which states the upper bound. Afterwards, we prove in Section 4 that the upper bound is tight (Theorem 2).

We begin with a basic definition that we will use throughout this section.

Definition. A knot or its location is *unique* if the knot’s input coordinate is different from that of all other knots in the neural network.

To set the base case for the induction, we first consider the neural network with $l = 1$ layer and n_1 neurons in that layer. Using the notation of Lemma 1, we make the simple observation that in a one-hidden-layer neural network, each neuron contributes exactly one knot to the model at $x_j = -b_{1j}/w_{1j}$. If the input biases b_{1j} and input weights w_{1j} are selected such that the knot locations x_j are unique, then the neural network has exactly n_1 knots.

To consider the inductive step, recall from Section 2.2 that every application of a rectified linear unit can preserve, destroy, or create new knots. For the purposes of constructing an upper bound, we can make the stronger statement that with the proper choice of weights and biases, every knot can be preserved in every hidden layer. Explicitly, the knots in the affine transformation $\mathbf{w}_{ik} \cdot \mathbf{v}_{i-1}(x) + b_{ik}$ of layer $i - 1$ outputs can be preserved in $\sigma(\mathbf{w}_{ik} \cdot \mathbf{v}_{i-1}(x) + b_{ik})$, the output of neuron k in layer i . The most naive way to do so is to set the biases b_{ik} so high that $\mathbf{w}_{ik} \cdot \mathbf{v}_{i-1}(x_j) + b_{ik} > 0$ for all knots x_j ; see Figure 5(a). The disadvantage of this method is that the rectified linear unit does not create any new knots. A better but still very simple alternative is to have two neurons in layer i employ identical or similar weights \mathbf{w}_{ik} and biases b_{ik} , but with flipped signs. This way, as shown in Figure 5(b), one neuron would preserve some subset of the knots of $\mathbf{w}_{ik} \cdot \mathbf{v}_{i-1}(x) + b_{ik}$, and the other neuron would preserve the complement. With this design, each rectified linear unit is able to create the maximum possible number of knots as follows.

Since each affine transformation $\mathbf{w}_{ik} \cdot \mathbf{v}_{i-1}(x) + b_{ik}$ is a linear spline, each line segment between adjacent knots can have at most one root. If $\mathbf{w}_{ik} \cdot \mathbf{v}_{i-1}(x) + b_{ik}$ has m_{i-1} knots, then these connections can cumulatively have at most $m_{i-1} - 1$ roots. Additionally, there may exist one root between $x = -\infty$ and the knot x_1 closest to $-\infty$, and another root between the knot $x_{m_{i-1}}$ closest to ∞ and $x = \infty$. In total, $\mathbf{w}_{ik} \cdot \mathbf{v}_{i-1}(x) + b_{ik}$ can have at most $m_{i-1} + 1$ roots. Hence, the output $\sigma(\mathbf{w}_{ik} \cdot \mathbf{v}_{i-1}(x) + b_{ik})$ of neuron k in layer i can create at most $m_{i-1} + 1$ knots, with the equality

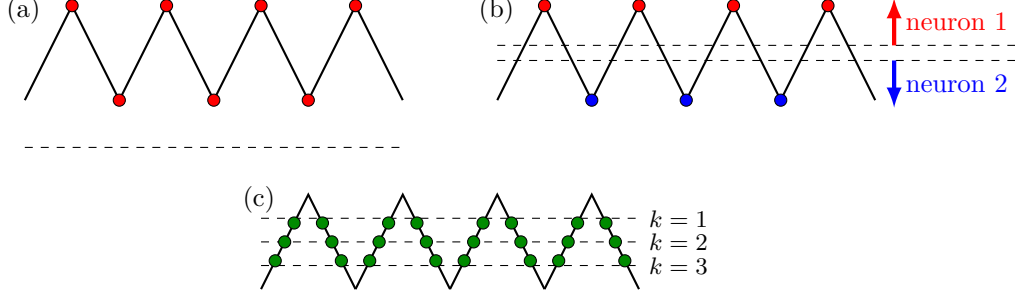


Figure 5: Schematics for preserving and creating knots in neuron k of layer i . (a) All knots x_j in $\mathbf{w}_{ik} \cdot \mathbf{v}_{i-1}(x) + b_{ik}$ (red) can be preserved in $\sigma(\mathbf{w}_{ik} \cdot \mathbf{v}_{i-1}(x) + b_{ik})$ by setting b_{ik} sufficiently high so that $\mathbf{w}_{ik} \cdot \mathbf{v}_{i-1}(x_j) + b_{ik}$ is greater than zero (dashed line) for all k . (b) Alternatively, two neurons (red and blue) can assign similar weights and biases with opposite signs to preserve knots on both sides of zero. (c) If $\mathbf{w}_{ik} \cdot \mathbf{v}_{i-1}(x) + b_{ik}$ is a sawtooth wave with m_{i-1} knots, then each neuron k in layer i can uniquely create $m_{i-1} + 1$ new knots. An example is shown for $k = 1, 2, 3$.

being met with a sawtooth wave. Furthermore, each neuron k can adjust b_{ik} so as to create $m_{i-1} + 1$ knots uniquely. This construction is demonstrated in Figure 5(c).

Having shown that all knots can be preserved in every layer, and having computed the maximum number of knots that each neuron can create, the upper bound (Theorem 1) can be formally derived. Note that we have not yet shown that all knots can always be preserved at the same time that every neuron in every layer creates the maximum possible number of knots. We first prove the upper bound as follows, and demonstrate the tightness of the bound by construction later in Section 4.

Proof of Theorem 1. For $l = 1$, the neural network can have up to one knot per neuron, as previously stated. That is, $m_1 \leq n_1$, which is equivalent to (1).

For $l > 1$, let us once again denote the number of knots in the affine transformation of layer i outputs by m_i . In layer i , each neuron $j = 1, \dots, n_i$ can preserve at most all m_{i-1} knots from the previous layer, and can also create at most $m_{i-1} + 1$ knots uniquely. Therefore, the upper bound on m_i is

$$m_i \leq m_{i-1} + n_i(m_{i-1} + 1) \quad (13a)$$

$$= (n_i + 1)m_{i-1} + n_i. \quad (13b)$$

Setting $i = l + 1$ in (13b), we have that $m_{l+1} \leq (n_{l+1} + 1)m_l + n_{l+1}$. Supposing that (1) is true, we find that

$$m_{l+1} \leq (n_{l+1} + 1) \sum_{i=1}^l n_i \prod_{j=i+1}^l (n_j + 1) + n_{l+1} \quad (14a)$$

$$= \sum_{i=1}^l n_i \prod_{j=i+1}^{l+1} (n_j + 1) + n_{l+1} \quad (14b)$$

$$= \sum_{i=1}^{l+1} n_i \prod_{j=i+1}^{l+1} (n_j + 1). \quad (14c)$$

Hence, if (1) holds for l , then it also holds for $l + 1$, and the induction is complete. \square

Remark. The dimension p of the output space does not affect the upper bound on the number of knots in the neural network; see Lemma 2 of Pascanu et al. (2014). The output layer is simply an

affine transformation, and does not contain any rectified linear units. Therefore, all knots that are outputted from the final hidden layer \mathbf{v}_l can be preserved. Additionally, some knots may possibly be destroyed in the degenerate case where \mathbf{v}_l has discontinuities in its first derivative, but $\mathbf{w}_{l+1,k} \cdot \mathbf{v}_l$ does not for all $k = 1, \dots, p$. Either way, no new knots can be created in the output layer.

Remark. In most applications of neural networks, $n_1 \geq \dots \geq n_l$, where n_l is notably larger than unity. In this case, the upper bound (1) is dominated by the $i = 1$ summand, and the upper bound is approximately

$$\prod_{i=1}^l n_i. \quad (15)$$

If we further assume that

$$n := n_1 = \dots = n_l \quad (16)$$

(which is sometimes useful for analytical purposes but less commonly employed in practice), then the upper bound further reduces to n^l . This approximate upper bound is consistent with the tight asymptotic upper bound $\mathcal{O}(n^{ql})$ given by Raghu et al. (2016), where we have used the input dimension $q = 1$.

Remark. The number of scalar parameters in the weights and biases of a deep $\mathbb{R}^q \rightarrow \mathbb{R}^p$ network (4) is

$$(q+1)n_1 + \sum_{i=1}^{l-1} (n_i + 1)n_{i+1} + (n_l + 1)p. \quad (17)$$

If we assume (16) once again, then for $q = 1$, the number of parameters is $2n + (n+1)(n(l-1)+p) \approx (p+2)n + (l-1)n^2$. This number is typically far smaller than n^l for $l \geq 3$. Thus, deep networks can possibly create a large number of knots with a comparatively small number of parameters. This feature plays a key role in the expressive power of deep neural networks. It has been suggested that although shallow networks can create models identical to deep networks via the universal approximation theorem, they may require many more parameters to do so; see Lin and Tegmark (2016) and the references within.

4 Tightness of the upper bound

Next, we show that the upper bound (1) is tight if there is a sufficient number of neurons in each layer, which will almost certainly be satisfied in practical applications. This demonstration proceeds by construction. In Lemma 2, we first review the trivial case where the neural network has $l = 1$ layer. We then show in Lemma 3 that the affine transformation of the first hidden layer outputs can be made into a sawtooth wave. Then, we show in Lemma 4 that subsequent hidden layers can turn sawtooth wave inputs into sawtooth wave outputs with the maximum number of knots. Finally, we reaffirm that all knots from a previous layer can be preserved in the application of a new layer, while creating the maximum number of knots.

Lemma 2. *The upper bound (1) is tight for single-hidden-layer neural networks.*

Proof. Equation (1) reduces to $m_1 \leq n_1$ for $l = 1$. As previously stated, the equality is obtained simply by choosing b_{1j} and w_{1j} in (5) such that $x_j = -b_{1j}/w_{1j}$ is unique for each $j = 1, \dots, n_1$. \square

Lemma 3. *If the first hidden layer has $n_1 \geq 3$ neurons, then there exist weights w_{1j}, w_{2kj} and biases b_{1j}, b_{2k} such that the input*

$$\sum_{j=1}^{n_1} w_{2kj} \sigma(w_{1j}x + b_{1j}) + b_{2k} \quad (18)$$

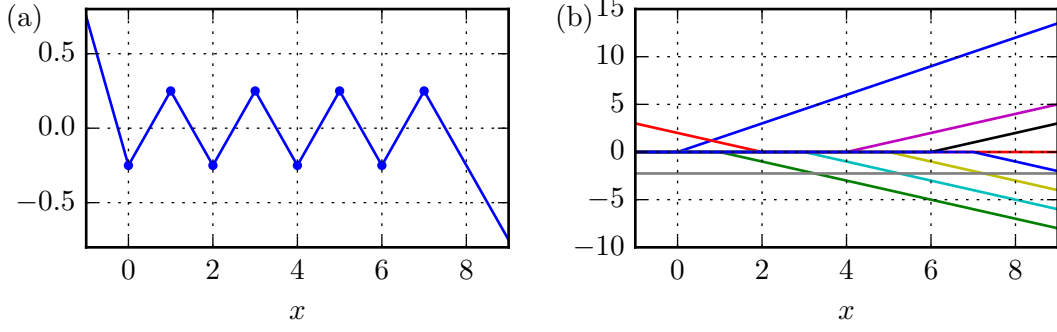


Figure 6: (a) The affine transformation (18) of first-hidden-layer outputs using the parameters in (19) with $n = 8$ and $b_{2k} = -9/4$. (b) The rectified linear unit summands in (a), with each summand in a different non-gray color (see (18)), and the bias b_{2k} in gray.

to the rectified linear unit in neuron k of layer 2 is a sawtooth wave.

Proof. One way to construct such a sawtooth wave is to select

$$w_{1j} = \begin{cases} -1 & | \quad j = 3 \\ 1 & | \quad j \neq 3 \end{cases} \quad (19a)$$

$$b_{1j} = \begin{cases} j - 1 & | \quad j = 3 \\ -j + 1 & | \quad j \neq 3 \end{cases} \quad (19b)$$

$$w_{2kj} = \begin{cases} \frac{3}{2} & | \quad j = 1 \\ -1 & | \quad j \text{ even} \\ 1 & | \quad j > 1 \text{ and } j \text{ odd} \end{cases}, \quad (19c)$$

with b_{2k} arbitrary. This is more apparent if we apply Lemma 1 and write (18) as

$$\sum_{j=1}^{n_1} s_{kj} \sigma(x - x_j) + c_{1k}x + c_{0k}, \quad (20)$$

where

$$x_j = j - 1, \quad s_{kj} = w_{2kj}, \quad c_{1k} = -1, \quad c_{0k} = b_{2k} + 2. \quad (21)$$

That is, the knots are evenly spaced, the initial slope from $x = -\infty$ to the first knot $x_1 = 0$ is $c_{1k} = -1$, and the slopes of the subsequent segments between knots are obtained by cumulatively adding s_{kj} . Thus, the slopes in successive linear pieces of the spline are

$$\left\{ c_{1k} + \sum_{j=1}^r s_{kj} \right\}_{r=0}^{n_1} = \left\{ -1, \frac{1}{2}, -\frac{1}{2}, \frac{1}{2}, -\frac{1}{2}, \dots \right\}, \quad (22)$$

which generates a sawtooth wave. See Figure 6 for an example. \square

Lemma 4. Suppose layer $i \geq 2$ has $n_i \geq 3$ neurons, and there exist weights $\alpha_{ij} \in \mathbb{R}$ for $j = 1, \dots, n_{i-1}$ such that

$$g_i(x) := \sum_{j=1}^{n_{i-1}} \alpha_{ij} v_{i-1,j}(x) \quad (23)$$

(which is an input to a layer i rectified linear unit, up to a bias) is a sawtooth wave with m_{i-1} knots. Then there exist weights $w_{ikj}, \alpha_{i+1,k} \in \mathbb{R}$ and biases $b_{ik} \in \mathbb{R}$ for $j = 1, \dots, n_{i-1}$ and $k = 1, \dots, n_i$ such that given

$$v_{ik}(x) := \sigma \left(\sum_{j=1}^{n_{i-1}} w_{ikj} v_{i-1,j}(x) + b_{ik} \right), \quad (24)$$

the function

$$g_{i+1}(x) := \sum_{k=1}^{n_i} \alpha_{i+1,k} v_{ik}(x) \quad (25)$$

is a sawtooth wave with the maximal number of knots

$$m_i = m_{i-1} + n_i(m_{i-1} + 1) \quad (26)$$

(cf. (13a)).

Proof. Suppose that—excluding the sections of $g_i(x)$ between $x = -\infty$ and the first knot x_1 , and between the last knot $x_{m_{i-1}}$ and $x = \infty$ —the minimum and maximum of the oscillation in $g_i(x)$ are respectively g_{\min} and g_{\max} . For convenience, let us rescale $g_i(x)$ such that the minimum and maximum are respectively 0 and 1; we define

$$\hat{g}_i(x) := \frac{g_i(x) - g_{\min}}{g_{\max} - g_{\min}}. \quad (27)$$

The central idea behind the construction is to select the weights and biases so that every line segment of the oscillation between $\hat{g}_i = 0$ and 1 is transformed into a sawtooth wave with n_i knots.

One method to achieve this is to construct the wave

$$\begin{aligned} g_{i+1}(x) = & \frac{3}{2} \sigma \left(\hat{g}_i(x) - \frac{1}{2n_i + 1} \right) - \sigma \left(\hat{g}_i(x) - \frac{3}{2n_i + 1} \right) \\ & + \sigma \left(-\hat{g}_i(x) + \frac{5}{2n_i + 1} \right) + \sum_{k=4}^{n_i} (-1)^{k+1} \sigma \left(\hat{g}_i(x) - \frac{2k-1}{2n_i + 1} \right). \end{aligned} \quad (28)$$

This construction has a natural equivalence with (19), with \hat{g}_i used in place of x . Interpreting \hat{g}_i as the independent variable and setting

$$\alpha_{i+1,k} := \begin{cases} \frac{3}{2} & | \quad k = 1 \\ -1 & | \quad k \text{ even} \\ 1 & | \quad k > 1 \text{ and } k \text{ odd} \end{cases}, \quad \gamma_k := \frac{2k-1}{2n_i + 1}, \quad (29)$$

we employ (9) to find that (28) is equivalent to

$$g_{i+1} = \sum_{k=1}^{n_i} \alpha_{i+1,k} \sigma(\hat{g}_i - \gamma_k) - \hat{g}_i + \frac{5}{2n_i + 1}. \quad (30)$$

Thus, as \hat{g}_i increases from 0 to 1, the slope of g_{i+1} with respect to \hat{g}_i in consecutive segments is

$$\left\{ -1 + \sum_{k=1}^r \alpha_{i+1,k} \right\}_{r=0}^{n_i} = \left\{ -1, \frac{1}{2}, -\frac{1}{2}, \frac{1}{2}, -\frac{1}{2}, \dots \right\}, \quad (31)$$

(cf. (22) and see Figure 7). Hence, for every line segment of $\hat{g}_i(x)$ between consecutive knots,

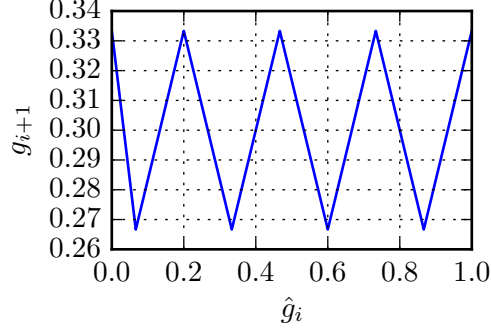


Figure 7: The wave (28) with $n_i = 7$.

$g_{i+1}(x)$ is a sawtooth wave with n_i knots.

Referring back to Section 3, we recall that the maximum number of knots (26) is achieved if every knot in $\hat{g}_i(x)$ is retained, and each of the n_i neurons uniquely creates $m_{i-1} + 1$ knots. We verify that these conditions are met. The quantity $\hat{g}_i - \gamma_k$ has a total of $m_{i-1} - 1$ roots between the m_{i-1} knots, plus one each between $x = -\infty$ and the first knot x_1 , and between the last knot $x_{m_{i-1}}$ and $x = \infty$. In total, each neuron creates $m_{i-1} + 1$ knots. Furthermore, each bias γ_k is unique, ensuring that the knots that are created by each of the n_i rectified linear units are also unique (see Figure 5(c)). Finally, since the operand to σ in the third summand in (28) contains $-\hat{g}_i(x)$ as opposed to $\hat{g}_i(x)$ in all other summands, both the lower and the upper knots of the sawtooth wave are preserved by the right-hand side of (28), as shown in Figure 5(b).

Note that for the induction to carry through successive layers, we must also verify that the local minima of (30) are all equal, as are the local maxima. This is easily confirmed, since the spacing in \hat{g}_i between consecutive knots (including endpoints) is

$$\{\gamma_1 - 0, \gamma_2 - \gamma_1, \dots, \gamma_{n_i} - \gamma_{n_i-1}, 1 - \gamma_{n_i}\} = \left\{ \frac{1}{2n_i + 1}, \frac{2}{2n_i + 1}, \dots, \frac{2}{2n_i + 1} \right\}. \quad (32)$$

Comparing this against the slopes (31), the vertical displacement between consecutive knots is simply

$$\left\{ -\frac{1}{2n_i + 1}, \frac{1}{2n_i + 1}, -\frac{1}{2n_i + 1}, \frac{1}{2n_i + 1}, \dots \right\}. \quad (33)$$

Finally, to complete the construction, we combine (23, 24, 27, 28) to find that one valid set of weights and biases is given by (29) and

$$w_{ikj} = \frac{\alpha_{ij}}{g_{\max} - g_{\min}} \cdot \begin{cases} -1 & | & k = 3 \\ 1 & | & k \neq 3 \end{cases} \quad (34a)$$

$$b_{ik} = - \left(\frac{g_{\min}}{g_{\max} - g_{\min}} + \frac{2k - 1}{2n_i + 1} \right) \cdot \begin{cases} -1 & | & k = 3 \\ 1 & | & k \neq 3 \end{cases}. \quad (34b)$$

□

With these lemmas in place, the tightness of the upper bound (Theorem 2) can now be proven.

Proof of Theorem 2. For $i = 1, \dots, l - 1$, the inductive and constructive proof is given quite simply by the combination of Lemmas 2–4. In the base case, Lemma 2 shows that (1) is tight for $l = 1$. Next, Lemma 3 shows that the affine transformation of the first hidden layer outputs—whether it

is for the output of a single-hidden-layer neural network, or for a second hidden layer in a deep network—can be made into a sawtooth wave. In light of Lemma 2, this sawtooth wave can be constructed with the maximal $m_1 = n_1$ knots.

Next, the induction step is given by Lemma 4. Namely, suppose that the affine transformation of the layer $i - 1$ outputs is a sawtooth wave with the maximal number of knots m_{i-1} . Then, it is possible to construct a sawtooth wave out of an affine transformation of the layer i outputs, such that the wave also has the maximal number of knots $m_i = m_{i-1} + n_i(m_{i-1} + 1)$. This induction step can be carried out sequentially from the second hidden layer $i = 2$ all the way to the penultimate hidden layer $i = l - 1$.

Finally, we note that the final hidden layer $i = l$ deserves special treatment because the output layer does not contain any rectified linear units. As a direct result, it is not actually necessary for the final hidden layer to output a sawtooth wave. Section 5 will later demonstrate this idea in an example. Instead, it is sufficient to have two neurons in the final hidden layer and still maintain the induction relation (13a). By referring back to Figure 5(b), we remind that two neurons can preserve all m_{i-1} knots from the penultimate layer, while each uniquely introducing $m_{i-1} + 1$ new knots with the application of the rectified linear unit. \square

In the constructive proofs of Lemmas 3 and 4, it is apparent that special consideration has been given to the third neuron in the respective series. This is also evident in Figure 6(b), which shows that the sawtooth wave in the affine transformation of the first hidden layer outputs can be constructed from all forward-facing rectified linear units, except for the third unit which faces backwards. To construct a sawtooth wave, it is in fact necessary to reverse the orientation of neuron j for some $j \geq 3$. Since a maximally high-wavenumber wave must be input into every rectified linear unit to meet the upper bound, an additional result is the following corollary, which is essentially the inverse of Theorem 2. We remark that the conditions of this corollary may not be seen in practice, but we nevertheless state this result for completeness.

Corollary 1. *For deep neural networks with $l \geq 2$ layers, the upper bound in (1) is not tight if $n_i < 3$ for any $i = 1, \dots, l - 1$, or if $n_l = 1$.*

Proof. For the upper bound to be met with $l \geq 2$, the affine transformations of the outputs of hidden layers $i = 1, \dots, l - 1$ must have alternating slopes—i.e., between positive and negative—through all linear pieces. Only then can each rectified linear unit in layer $i + 1$ create the maximal $m_i + 1$ unique knots. This condition can be analyzed separately for $i = 1$ and $i > 1$.

For $i = 1$, the individual rectified linear units of the first hidden layer must be linearly combined to construct a sawtooth wave; see Lemma 3. Such an arrangement is not possible in the (rather unorthodox) case of $n_1 = 1$ or 2. The case where $n_1 = 1$ is trivial: the function $\sigma(w_{11}x + b_{11})$ clearly cannot have both a negative and a positive slope for a given choice of w_{11} and b_{11} . The case where $n_1 = 2$ is slightly less obvious. Suppose, without loss of generality, that we wish to construct a linear combination

$$g_2(x) = \sum_{j=1}^2 w_{2j}\sigma(w_{1j}x + b_{1j}) \quad (35)$$

of two neural outputs in the first hidden layer that slopes down, then up, and finally down again: $\searrow \swarrow$. The left and right extremes of this shape requires that one neuron be oriented toward quadrant II (\searrow) and the second neuron be oriented toward quadrant IV (\swarrow). There does not exist a way to sum these two rectified linear units and obtain the positive slope in the middle segment of the linear combination. Therefore, the upper bound (1) cannot be achieved if $n_1 < 3$.

For $i = 2, \dots, l - 1$, hidden layer i must be able to transform a sawtooth wave with m_{i-1} knots into another sawtooth wave with $m_{i-1} + n_i(m_{i-1} + 1)$ knots. Consider a single line segment in the

linear combination of layer $i - 1$ outputs. Using the notation of Lemma 4, if the output of this segment has a minimum $g_i = g_{\min}$ and maximum $g_i = g_{\max}$, then we require some choice of w_{ij} , $w_{i+1,j}$ and b_{ij} such that the derivative of

$$g_{i+1} = \sum_{j=1}^{n_i} w_{i+1,j} \sigma(w_{ij} g_i + b_{ij}) \quad (36)$$

with respect to g_i contains n_i sign changes as g_i increases from g_{\min} to the next instance of g_{\max} . Using the same argument as the previous paragraph for $i = 1$, but using the input g_i in place of x , such an arrangement is impossible if $n_i = 1$ or 2 .

Finally, we make the observation that in the unusual case that $n_l = 1$, it is impossible for that single final-hidden-layer neuron both to preserve all m_{l-1} knots from the penultimate layer, while also introducing $m_{l-1} + 1$ knots. If $m_{l-1} + 1$ knots were introduced by drawing a bias through the sawtooth wave from layer $l - 1$, then half of the m_{l-1} knots (rounded up or down, if m_{l-1} is odd) from the previous layer would be discarded. Alternatively, if the single neuron preserved all m_{l-1} knots from the previous layer, then it would not be able to create new knots, as required by the upper bound. \square

5 Example construction of tight upper bound

In this section, we demonstrate a construction of an $\mathbb{R} \rightarrow \mathbb{R}^p$ neural network with a number of knots exactly equal to the upper bound. For the sake of keeping the neural network size manageable, we intentionally use a small number of neurons. We choose to have $l = 3$ hidden layers, with $n_1 = 6$ neurons in the first layer, $n_2 = 3$ neurons in the second layer, and $n_3 = 2$ neurons in the third layer. We will employ $p = 2$ in this example, though as Section 3 shows, the output dimension is actually irrelevant to the number of knots in the neural network.

Using these values in (1), we find that the upper bound on the number of knots is $m_1 = 6$ in the first layer outputs, $m_2 = 27$ in the second layer outputs, and $m_3 = 83$ in the third layer and final outputs. Since n_1 , n_2 , and n_3 satisfy the criteria in Theorem 2, these bounds are tight, and we can use the constructions in Section 4 to define a neural network with these numbers of knots.

The example neural network is given by the equations

$$v_{1k} = \sigma(w_{1k}x + b_{1k}), \quad k = 1, \dots, n_1 \quad (37a)$$

$$v_{2k} = \sigma \left(\sum_{j=1}^{n_1} w_{2kj} v_{1j} + b_{2k} \right), \quad k = 1, \dots, n_2 \quad (37b)$$

$$v_{3k} = \sigma \left(\sum_{j=1}^{n_2} w_{3kj} v_{2j} + b_{3k} \right), \quad k = 1, \dots, n_3 \quad (37c)$$

$$y_k = \sum_{j=1}^{n_3} w_{4kj} v_{3j} + b_{4k}, \quad k = 1, \dots, p, \quad (37d)$$

where

$$w_{1k} = \begin{cases} -1 & | & k = 3 \\ 1 & | & k \neq 3 \end{cases}, \quad (38a)$$

$$b_{1k} = \begin{cases} k - 1 & | & k = 3 \\ -k + 1 & | & k \neq 3 \end{cases} \quad (38b)$$

for $k = 1, \dots, n_1$,

$$w_{2kj} = 2w_{1k} \cdot \begin{cases} \frac{3}{2} & | & j = 1 \\ -1 & | & j \text{ even} \\ 1 & | & j > 1 \text{ and } j \text{ odd} \end{cases}, \quad (39a)$$

$$b_{2k} = \left(-4 - \frac{2k-1}{2n_2+1}\right) w_{1k} \quad (39b)$$

for $k = 1, \dots, n_2$,

$$w_{3kj} = 7(-1)^{k-1} \cdot \begin{cases} \frac{3}{2} & | & j = 1 \\ -1 & | & j \text{ even} \\ 1 & | & j > 1 \text{ and } j \text{ odd} \end{cases}, \quad (40a)$$

$$b_{3k} = (-1)^{k-1} \left(-4 - \frac{k}{n_3+1}\right) \quad (40b)$$

for $k = 1, \dots, n_3$, and

$$w_{4kj} = (-1)^{j+k}, \quad (41a)$$

$$b_{4k} = k - 1 \quad (41b)$$

for $k = 1, \dots, p$.

The hidden layer and model outputs for this example are shown in Figure 8. The interpretation of the above weights and biases proceeds as follows. In the first hidden layer, w_{1k} and b_{1k} (38), as well as the dependence of w_{2kj} (39a) on j , are copied directly from the construction for a sawtooth wave (19) in Lemma 3. Thus, they create knots at $x = 0, \dots, n_1 - 1$, and the rectified linear units are oriented as in Figure 6(b). The factor of 2 in (39a) is added for convenience to make the sawtooth span a range of 1 instead of $1/2$. The sawtooth wave

$$g_2(x) = \sum_{j=1}^{n_1} w_{21j} v_{1j}(x) \quad (42)$$

that is used in neuron $k = 1$ of layer $i = 2$ is shown in Figure 8(a).

From this figure, we observe that the range of the sawtooth wave, excluding the end parts with $g_2 \rightarrow \pm\infty$, is $[4, 5]$. Following (34), we flip the signs of w_{2kj} and b_{2k} (39) for $k = 3$. Furthermore, we set b_{2k} according to (34b), so that each neuron offsets g_2 by the proper amount to construct $m_1 + 1$ unique knots, which can then be rearranged into a new sawtooth wave. In addition, we set the dependence of w_{3kj} (40a) on j to match the construction in (28). As shown in Figure 8(b), this choice of parameters produces the sawtooth wave

$$g_3(x) = \sum_{j=1}^{n_2} w_{31j} v_{2j}(x) \quad (43)$$

that is used in neuron $k = 1$ of layer $i = 3$. We may observe from this figure that this second layer output retains all the knots from the first layer output (Figure 8(a)), and it also creates the maximal n_2 knots between all the knots of the first layer output, as well as in $(-\infty, 0)$ and $(n_1 - 1, \infty)$.

Moving forward, the construction of the third hidden layer in this example proceeds differently. As stated in Theorem 2, the final hidden layer $i = 3$ only needs to have $n_3 = 2$ neurons to meet

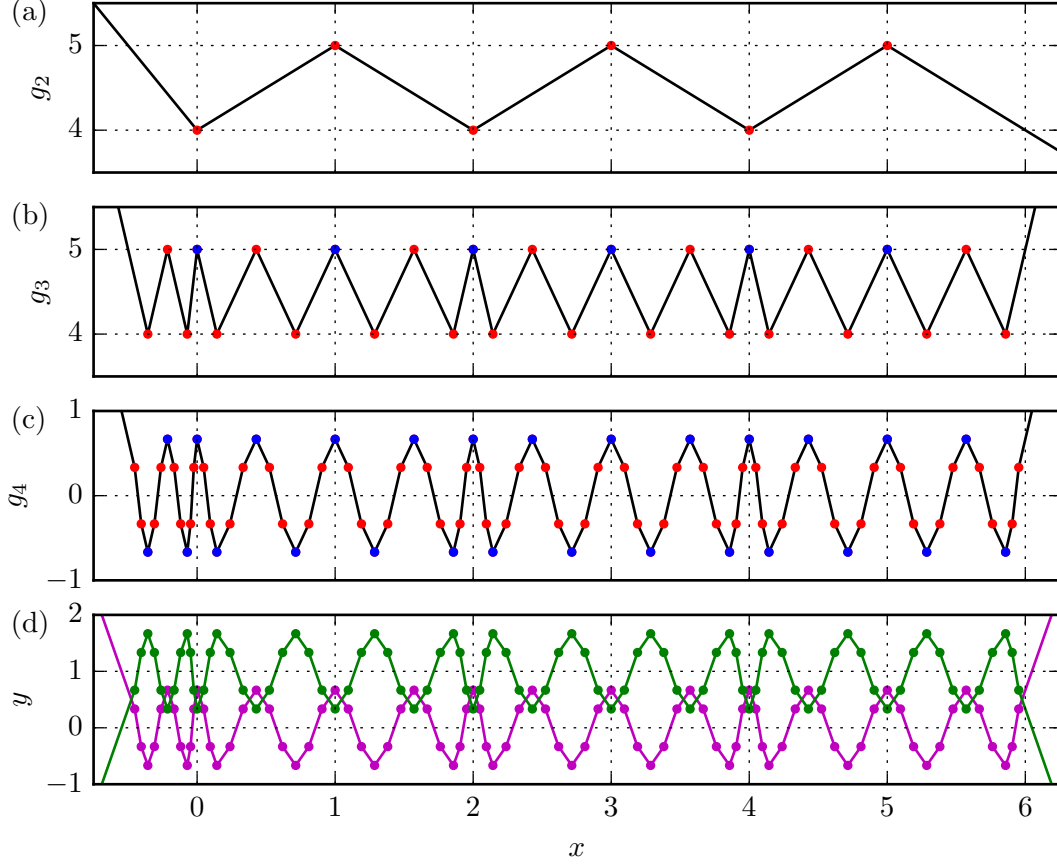


Figure 8: The neural network given by (37–41), as an example of a model that meets the upper bound (1) on the number of knots. The sawtooth waves $\sum_{j=1}^{n_i} w_{i+1,1,j} v_{ij}$ are constructed by linearly combining the outputs of hidden layer (a) $i = 1$, with six knots; (b) $i = 2$, with 27 knots; and (c) $i = 3$, with 83 knots. Knots retained from the previous layer are shown in blue, and knots created in the current layer are shown in red. (d) The outputs y_1 (magenta) and y_2 (green), with 83 knots.

the tight upper bound, since there are no further rectified linear units and the sawtooth waveform is therefore no longer required. By following the strategy shown in Figure 5(c), we pick w_{3kj} and b_{3k} (40) to have opposite signs between $k = 1$ and 2. Furthermore, we note that the sawtooth in Figure 8(b) has a range of $[4, 5]$, so we pick b_{3k} to be two different values for $k = 1$ and 2 within the range $(-5, -4)$. That way, as shown in Figure 5(b), the $k = 1$ neuron retains the upper knots of Figure 8(b), while the $k = 2$ neuron retains the lower ones. Furthermore, each of the two neurons produces one new knot in the $m_2 + 1$ regions of \mathbb{R} divided by the knots of Figure 8(b). The factor of seven in (40a) is arbitrary.

Finally, the choice of the output weights w_{4kj} and biases b_{4k} (41) is also arbitrary, since the output layer does not contain rectified linear units and cannot destroy or create knots. The sawtooth wave

$$g_4(x) = \sum_{j=1}^{n_3} w_{41j} v_{3j}(x) \quad (44)$$

that makes up the output y_1 is shown in Figure 8(c). The neural network outputs (37d), with the maximal $m_3 = 83$ knots, are shown in Figure 8(d).

6 Conclusion

We have shown that deep, fully-connected, $\mathbb{R} \rightarrow \mathbb{R}^p$ neural networks with rectified linear unit activations are essentially linear splines. In Theorem 1, we derived an upper bound on the number of knots that such neural networks can have. The upper bound is given exactly by (1); to close approximation, this bound is $n_1 \cdots n_l$. We then showed in Theorem 2 that the upper bound is tight for the neural network widths that would be encountered in practice. An example of a deep neural network exactly meeting this upper bound was described in Section 5.

It is clear from the setup of the upper bound that the imposed conditions are prohibitively restrictive. Most notably, it is common in practical applications to construct $\mathbb{R}^q \rightarrow \mathbb{R}^p$ neural networks where q may be on the order of 10^3 or even larger. As aforementioned, previous works have computed approximate or asymptotic bounds on the number of linear pieces in $\mathbb{R}^q \rightarrow \mathbb{R}^p$ neural networks (Montúfar et al. 2014; Pascanu et al. 2014; Raghu et al. 2016). Nevertheless, an exact upper bound—let alone a tight one—remains to be derived in this generic case.

In addition, there is little reason to believe that neural networks used in actual applications would contain a number of knots equal to or close to the upper bound presented here. The construction of the upper bound required that a sawtooth wave be constructed at every hidden layer except the final one. It is unlikely that such maximally high-wavenumber networks would be fitted to actual data, and the likelihood is even lower for large input dimensions q commonly used in practice.

Thus, the results of this paper can be interpreted as a theoretical “brick-wall” limit on neural network expressivity, which may be used as a guideline or check in designing actual neural networks. Two companion papers present more realistic scenarios. In the first (Chen et al. 2016a), we explore the number of knots in randomly weighted and biased neural networks. In the second (Chen et al. 2016b), we describe empirical results on the behavior of neural network training. Both of these scenarios are more representative of actual situations seen in practice. Not only is a random neural network more likely to represent an “average case” neural network rather than a “best case,” but also—as demonstrated in Chen et al. (2016a)—random neural networks are actually encountered in the early stages of training on data. In Chen et al. (2016a), we also describe open problems related to the expressivity of neural networks in greater detail. These papers are still largely analytical in nature, since the chief objective of our investigation is to close the gap between our understanding of neural network theory and applications.

Alden Walker is gratefully acknowledged for providing Figure 4 and for helpful conversations. Discussions with Anthony Gamst were also very fruitful, and led to the central ideas of the work presented in Chen et al. (2016a,b).

References

- K. K. Chen, A. C. Gamst, and A. K. Walker. Knots in random neural networks. In *Neural Information Processing Systems (NIPS), Workshop on Bayesian Deep Learning*, Barcelona, Spain, 2016a. To be presented.
- K. K. Chen, A. C. Gamst, and A. K. Walker. The empirical size and risk of trained neural networks, 2016b. arXiv:1611.09444.
- G. Cybenko. Approximation by superpositions of a sigmoidal function. *Math. Control Signals Syst.*, 2(4):303–314, 1989.

- K. Hornik. Approximation capabilities of multilayer feedforward networks. *Neural Netw.*, 4(2): 251–257, 1991.
- K. Hornik, M. Stinchcombe, and H. White. Multilayer feedforward networks are universal approximators. *Neural Netw.*, 2:359–366, 1989.
- S. W. Knox. Machine learning: Topics and techniques, Edition 2.2, 2016.
- H. W. Lin and M. Tegmark. Why does deep and cheap learning work so well?, 2016. arXiv:1608.08225v1.
- G. Montúfar, R. Pascanu, K. Cho, and Y. Bengio. On the number of linear regions of deep neural networks. In *Advances in Neural Information Processing Systems*, pages 2924–2932, 2014.
- R. Pascanu, G. Montúfar, and Y. Bengio. On the number of response regions of deep feedforward networks with piecewise linear activations, 2014. arXiv:1312.6098v5.
- M. Raghu, B. Poole, J. Kleinberg, S. Ganguli, and J. Sohl-Dickstein. On the expressive power of deep neural networks, 2016. arXiv:1606.05336v2.
- S. Sonoda and N. Murata. Neural network with unbounded activation functions is universal approximator. *Appl. Comput. Harmon. Anal.*, 2015. In press.